

MEDICAMENTOS

VIGENTES



JUAN JOSÉ RINCÓN MÉNDEZ - 2202018

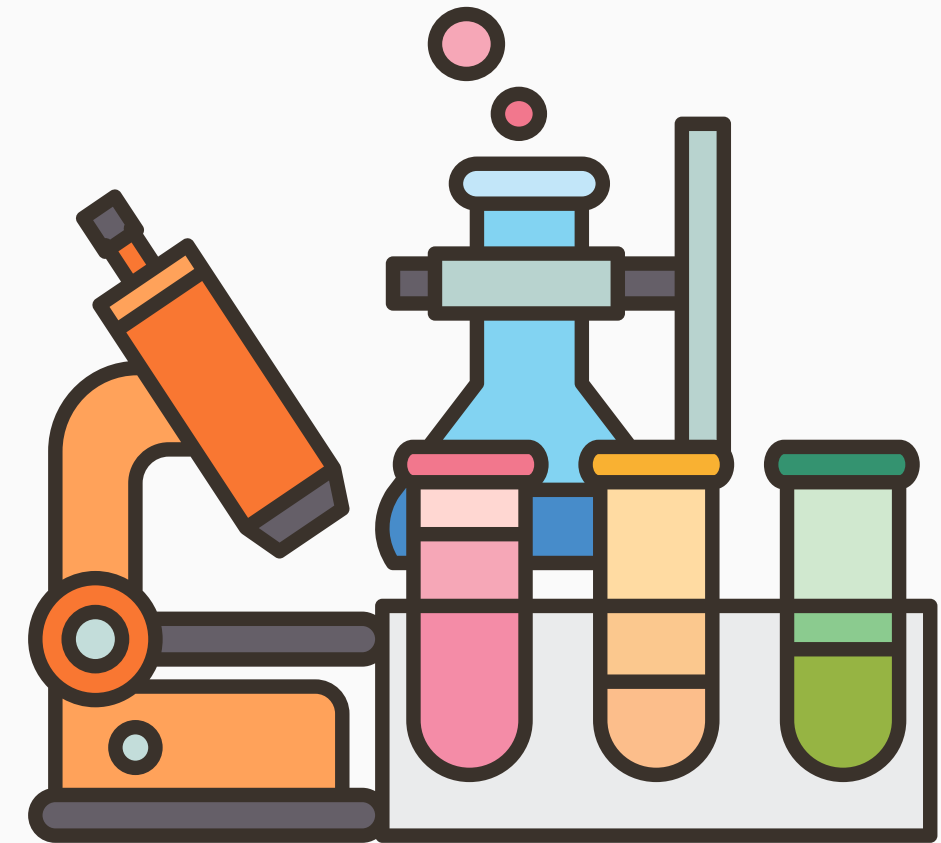
DYLAN FRANCISCO JIMENEZ SANDOVAL - 2202049

SONIA MARCELA GRANADOS MORENO - 2204250



DE QUE TRATA EL DATASET

Trata de medicamentos vigentes y lo que buscamos saber con que laboratorio un medicamento tiene mayor tiempo de duración para tener en cuenta a la hora de comprar.



Mounted at /content/drive

```
[ ] 1 import pandas as pd
    2 import numpy as np
    3 from datetime import datetime
    4 import datetime as dt
    5 import seaborn as sns
    6 import matplotlib.pyplot as plt
    7 from sklearn.tree import DecisionTreeClassifier
```

Cargar el dataset

```
[ ] 1 datos = pd.read_csv('/content/drive/MyDrive/medicamentos/medicamentos.csv')
    2
    3 datos = datos.replace(',', '.', regex=True)
    4
    5 datos.to_csv("dataset_modificado.csv", index=False)
    6
    7
```

```
[ ] 1 datos['expediente'] = datos['expediente'].astype(int)
    2 datos['producto'] = datos['producto'].astype(str)
```

Importamos las librerías que vamos a utilizar durante el desarrollo del proyecto y subimos el dataset.

Convertimos expediente y producto en datos int.

```
[ ] 1 datos['fechavencimiento'] = pd.to_datetime(datos['fechavencimiento'])
    2 datos['fechaexpedicion'] = pd.to_datetime(datos['fechaexpedicion'])
    3 datos['fechaactivo'] = pd.to_datetime(datos['fechaactivo'])
    4 datos['fechainactivo'] = pd.to_datetime(datos['fechainactivo'])
    5
    6 default_date = pd.Timestamp('00:00:00')
    7 datos['fechavencimiento'] = datos['fechavencimiento'].fillna(default_date)
    8 datos['fechaexpedicion'] = datos['fechaexpedicion'].fillna(default_date)
    9
```

```
[ ] 1 del datos['cantidadcum']
    2 del datos['expedientecum' ]
    3 del datos['concentracion' ]
    4 del datos['descripcioncomercial']
    5 del datos['unidad' ]
    6 del datos['atc' ]
    7 del datos['descripcionatc' ]
    8 del datos['viaadministracion']
    9 del datos['unidadmedida' ]
   10 del datos['unidadreferencia' ]
   11 del datos['formafarmaceutica']
   12 del datos['muestramedica' ]
   13 del datos['cantidad' ]
   14 del datos['IUM' ]
   15 del datos['consecutivocum']
```

Se busca eliminar algunas columnas del Dataset propuesto, ya que estas generan interferencia al momento de realizar predicciones o de entrenar a la maquina.

Estas líneas eliminan las columnas con los nombres especificados del DataFrame.

```
[7] 1 # Eliminar filas duplicadas en el DataFrame
    2 datos = datos.drop_duplicates()
```

```
1 datos.head(50)
```

51	7028	BIODERM® CREMA	LABORATORIOS SIEGFRIED S.A.S.	INVIMA 2023M-005385-R3	19
52	7028	BIODERM® CREMA	LABORATORIOS SIEGFRIED S.A.S.	INVIMA 2023M-005385-R3	19
59	7028	BIODERM® CREMA	LABORATORIOS SIEGFRIED S.A.S.	INVIMA 2023M-005385-R3	19
60	7028	BIODERM® CREMA	LABORATORIOS SIEGFRIED S.A.S.	INVIMA 2023M-005385-R3	19
61	7028	BIODERM® CREMA	LABORATORIOS SIEGFRIED S.A.S.	INVIMA 2023M-005385-R3	19

Eliminamos las filas duplicadas y las mostramos.

Después de ejecutar estas líneas de código, las columnas 'fechaexpedicion' y 'fechavencimiento' del DataFrame contendrán valores numéricos que representan el timestamp de las fechas originales. Con esto logrando obtener el promedio que es la columna con la cual se busca graficar

```
2 datos['fechavencimiento'] = datos['fechavencimiento'].apply(lambda x: x.timestamp())
```

Promedio

```
[10] 1 PromedioFechas = datos[['fechaexpedicion', 'fechavencimiento']].mean(axis=1)
    2 datos['Promedio'] = PromedioFechas
    3
```

```
1 datos.head(50)
```

16	5581	VAXOM®CAPSULAS ADULTOS	OM PHARMA SA	INVIMA 2023MB-007838-R4	915
17	5581	BRONCHO-VAXOM®CAPSULAS ADULTOS	OM PHARMA SA	INVIMA 2023MB-007838-R4	915
20	5581	BRONCHO-VAXOM®CAPSULAS ADULTOS	OM PHARMA SA	INVIMA 2023MB-007838-R4	915
21	5581	BRONCHO-VAXOM®CAPSULAS ADULTOS	OM PHARMA SA	INVIMA 2023MB-007838-R4	915


```
✓ [12] 1 datos.shape
```

```
(38698, 15)
```

```
✓ [13] 1 datos.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 38698 entries, 0 to 114705
Data columns (total 15 columns):
#   Column                Non-Null Count  Dtype
---  -
0   expediente            38698 non-null  int64
1   producto              38698 non-null  object
2   titular               38698 non-null  object
3   registrosanitario     38698 non-null  object
4   fechaexpedicion       38698 non-null  float64
5   fechavencimiento      38698 non-null  float64
6   estadoregistro        38698 non-null  object
7   estadocum             38698 non-null  object
8   fechaactivo           38698 non-null  datetime64[ns]
9   fechainactivo         8033 non-null   datetime64[ns]
10  principioactivo       38698 non-null  object
11  nombrerol             38698 non-null  object
12  tiporol               38698 non-null  object
13  modalidad             38698 non-null  object
14  Promedio              38698 non-null  float64
dtypes: datetime64[ns](2), float64(3), int64(1), object(9)
memory usage: 4.7+ MB
```

Obtenemos los tipos de datos

```
✓ [14] 1 datos.dtypes
```

```
expediente            int64
producto              object
titular               object
registrosanitario     object
fechaexpedicion       float64
fechavencimiento      float64
estadoregistro        object
estadocum             object
fechaactivo           datetime64[ns]
fechainactivo         datetime64[ns]
principioactivo       object
nombrerol             object
tiporol               object
modalidad             object
Promedio              float64
dtype: object
```

```
✓ [15] 1 datos.describe()
```

	expediente	fechaexpedicion	fechavencimiento	Promedio
count	3.869800e+04	3.869800e+04	3.869800e+04	3.869800e+04
mean	1.708055e+07	1.263123e+09	1.785009e+09	1.524066e+09
std	7.078071e+06	2.499299e+08	4.801744e+07	1.265010e+08
min	3.521000e+03	6.005664e+08	1.673309e+09	1.170461e+09
25%	1.993513e+07	1.059091e+09	1.746230e+09	1.422576e+09
50%	2.001021e+07	1.269994e+09	1.784160e+09	1.532390e+09
75%	2.009082e+07	1.460592e+09	1.829520e+09	1.625778e+09

Obtenemos las dimensiones del dataframe, una descripción concisa incluyendo el tipo de datos de cada columna, la cantidad de valores no nulos y la cantidad total de memoria utilizada.

Mostramos los datos

```
[15] 1 datos.expediente.describe()
```

```
count    3.869800e+04  
mean     1.708055e+07  
std       7.078071e+06  
min       3.521000e+03  
25%      1.993513e+07  
50%      2.001021e+07  
75%      2.009082e+07  
max       2.024415e+07  
Name: expediente, dtype: float64
```

```
[17] 1 datos.columns
```

```
Index(['expediente', 'producto', 'titular', 'registrosanitario',  
      'fechaexpedicion', 'fechavencimiento', 'estadoregistro', 'estadocum',  
      'fechaactivo', 'fechainactivo', 'principioactivo', 'nombrerol',  
      'tiporol', 'modalidad', 'Promedio'],  
      dtype='object')
```

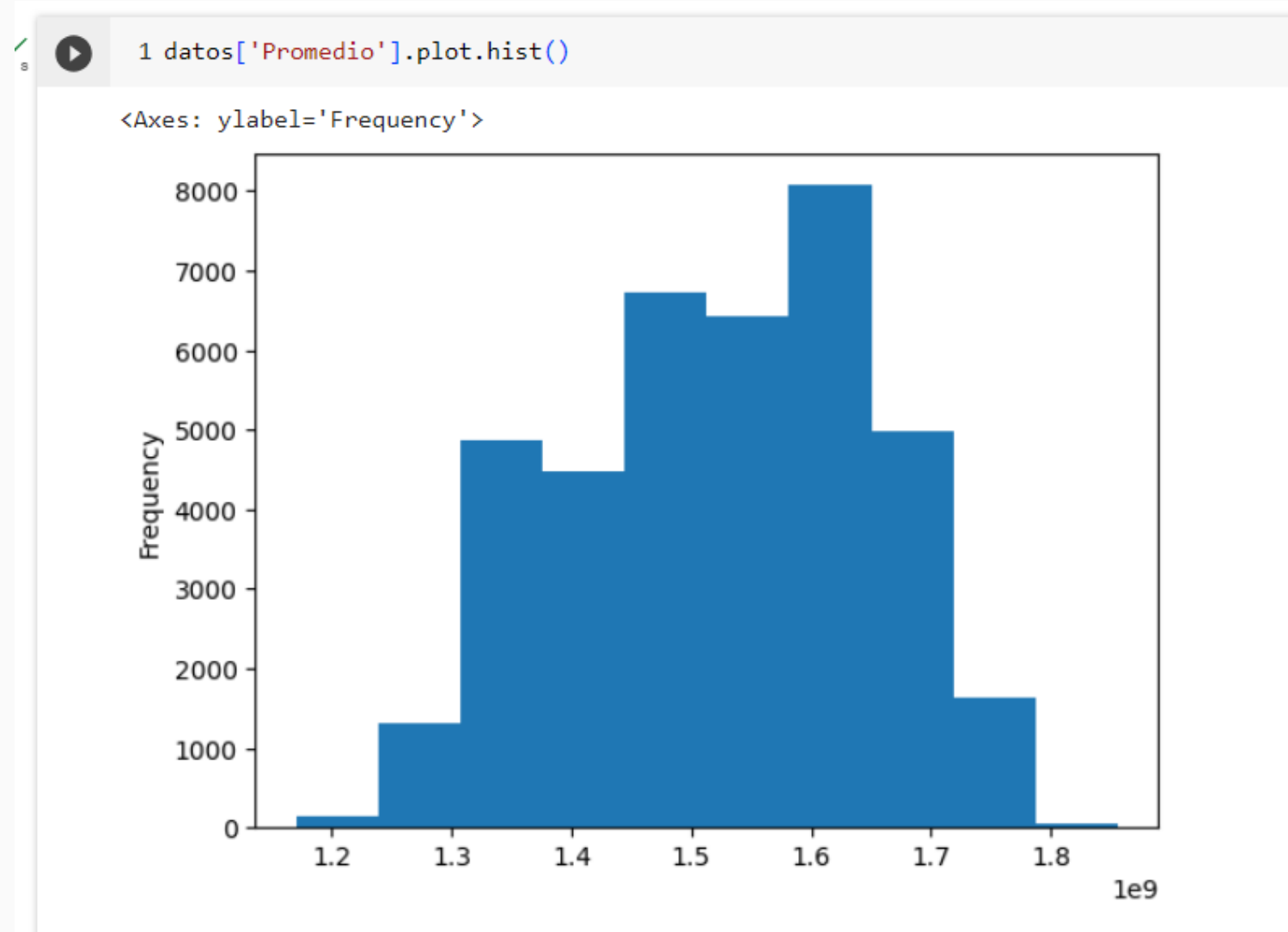
```
[18] 1 columnas = datos.columns  
2 columnas
```

```
Index(['expediente', 'producto', 'titular', 'registrosanitario',  
      'fechaexpedicion', 'fechavencimiento', 'estadoregistro', 'estadocum',  
      'fechaactivo', 'fechainactivo', 'principioactivo', 'nombrerol',  
      'tiporol', 'modalidad', 'Promedio'],  
      dtype='object')
```

```
1 datos.tail()
```

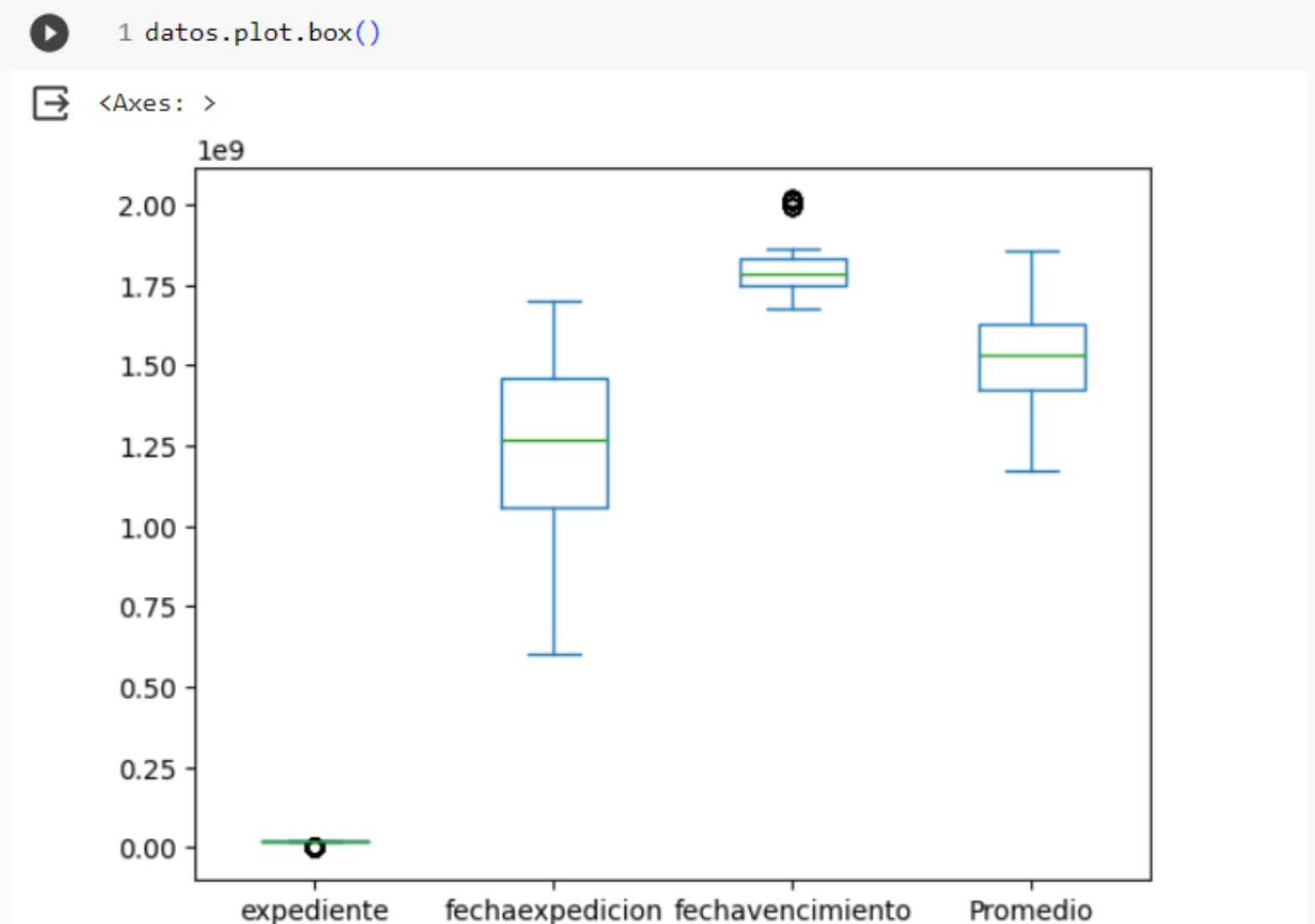
	expediente	producto	titular	registrosanitario	fechaexpedicion	fe
114701	20239915	MILACID® SACHET POR 10ML SABOR MENTA	BIOFLUIDOS & FARMA S.A.S.	INVIMA 2023M- 0020967	1.693872e+09	
114702	20239915	MILACID® SACHET POR 10ML SABOR MENTA	BIOFLUIDOS & FARMA S.A.S.	INVIMA 2023M- 0020967	1.693872e+09	
114703	20239915	MILACID® SACHET POR 10ML SABOR MENTA	BIOFLUIDOS & FARMA S.A.S.	INVIMA 2023M- 0020967	1.693872e+09	
114704	20243782	KEPINIA 25 MG	SALUS PHARMA LABS S.A.S.	INVIMA 2023M- 0020955	1.680653e+09	
114705	20244153	KEPINIA 100 MG	SALUS PHARMA LABS S.A.S.	INVIMA 2023M- 0020954	1.680653e+09	

Obtenemos un resumen estadístico de la columna expediente y las variables.



Realizamos un histograma, el cual mostrará la frecuencia de los diferentes valores en la columna 'Promedio'. Cada barra en el histograma representa un rango de valores, y la altura de la barra indica cuántas veces los valores en ese rango aparecen en la columna.

- Realizamos un diagrama de caja que representa el rango intercuartílico (IQR), que es la diferencia entre el tercer cuartil (Q3) y el primer cuartil (Q1). La línea en el medio de la caja es la mediana (Q2).
- Los "bigotes" se extienden hasta los valores que están dentro de 1.5 veces el IQR desde Q1 y Q3. Puntos fuera de estos límites se consideran valores atípicos y se representan como puntos individuales.

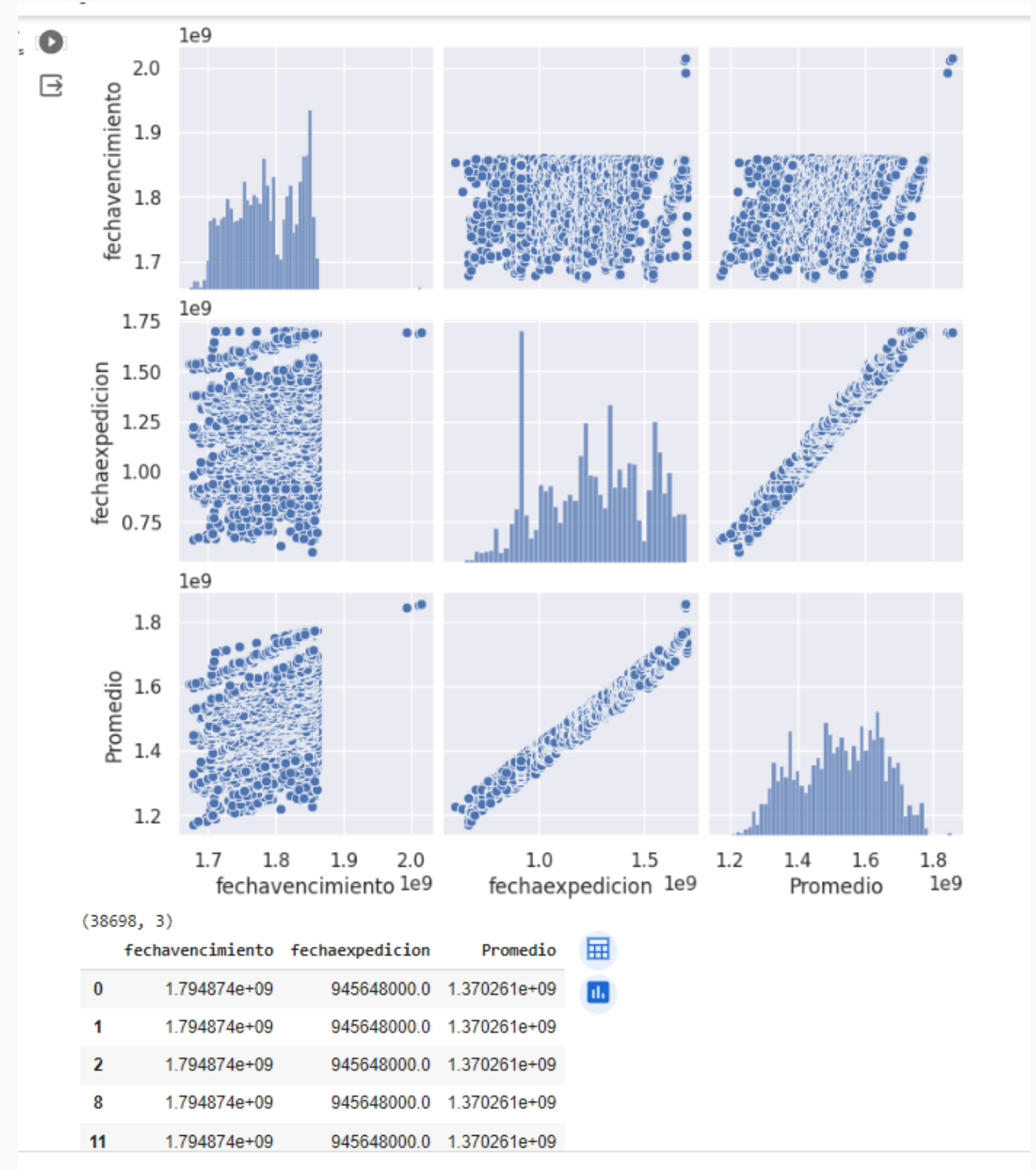



```

1 X = datos.values[:, 1:]
2 y = datos.values[:, [1]]
3
4 sns.set(rc={'figure.figsize': (15, 15)})
5
6 sns.pairplot(datos, diag_kind="hist")
7
8 plt.show()
9
10 print(X.shape, y.shape)
11
12 datos.head(5)
13
14 columns = ["fechavencimiento", "fechaexpedicion", "Promedio"]
15 X = datos[columns].values
16
17 sns.set(rc={'figure.figsize': (15, 15)})
18
19 sns.pairplot(datos[columns], diag_kind="hist") # kind="kde"
20
21 plt.show()
22
23 print(X.shape)
24 datos[columns].head(5)

```

Visualizamos la distribución y las relaciones entre las variables del DataFrame utilizando pair plots, además de seleccionar columnas específicas para un análisis más detallado.



```

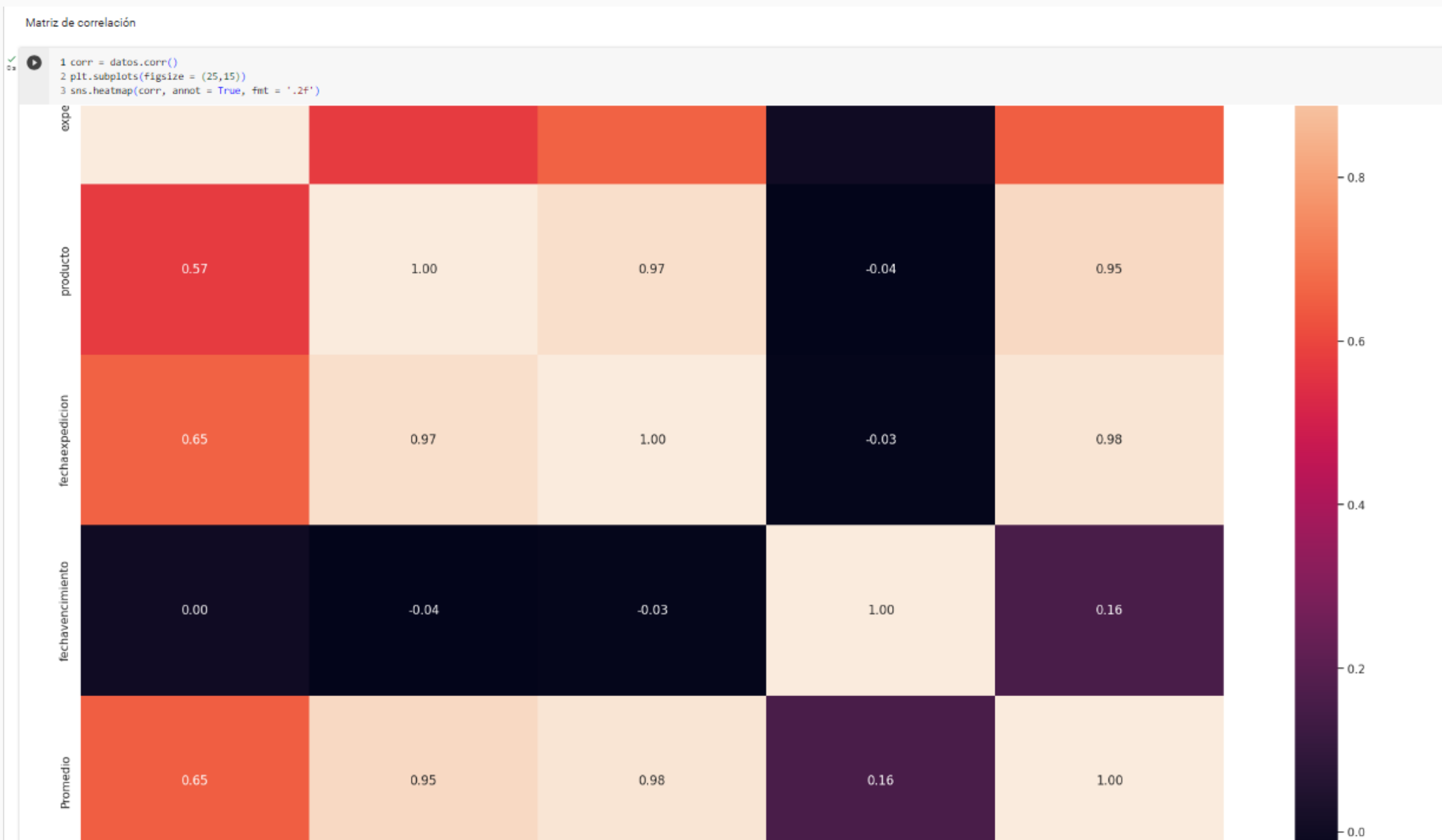
1 k=0
2 for i in datos.producto.unique():
3     datos.producto.replace(i, k, inplace = True)
4     k = k+1

```

```
1 datos.head(10)
```

	expediente	producto	titular	registrosanitario	fechaexpedicion	fechavenc
0	3521	0	PROCAPS S.A.	INVIMA 2021M-002103-R3	945648000.0	1.794
1	3521	0	PROCAPS S.A.	INVIMA 2021M-002103-R3	945648000.0	1.794
2	3521	0	PROCAPS S.A.	INVIMA 2021M-002103-R3	945648000.0	1.794
8	3521	0	PROCAPS S.A.	INVIMA 2021M-002103-R3	945648000.0	1.794
11	3521	0	PROCAPS S.A.	INVIMA 2021M-002103-R3	945648000.0	1.794
12	3521	0	PROCAPS S.A.	INVIMA 2021M-002103-R3	945648000.0	1.794
13	3521	0	PROCAPS S.A.	INVIMA 2021M-002103-R3	945648000.0	1.794
16	5581	1	OM PHARMA SA	INVIMA 2023MB-007838-R4	915148800.0	1.851

Este código proporciona una visualización del Dataset con los cambios realizados, donde se puede observar el cambio de fecha a un valor numerico.

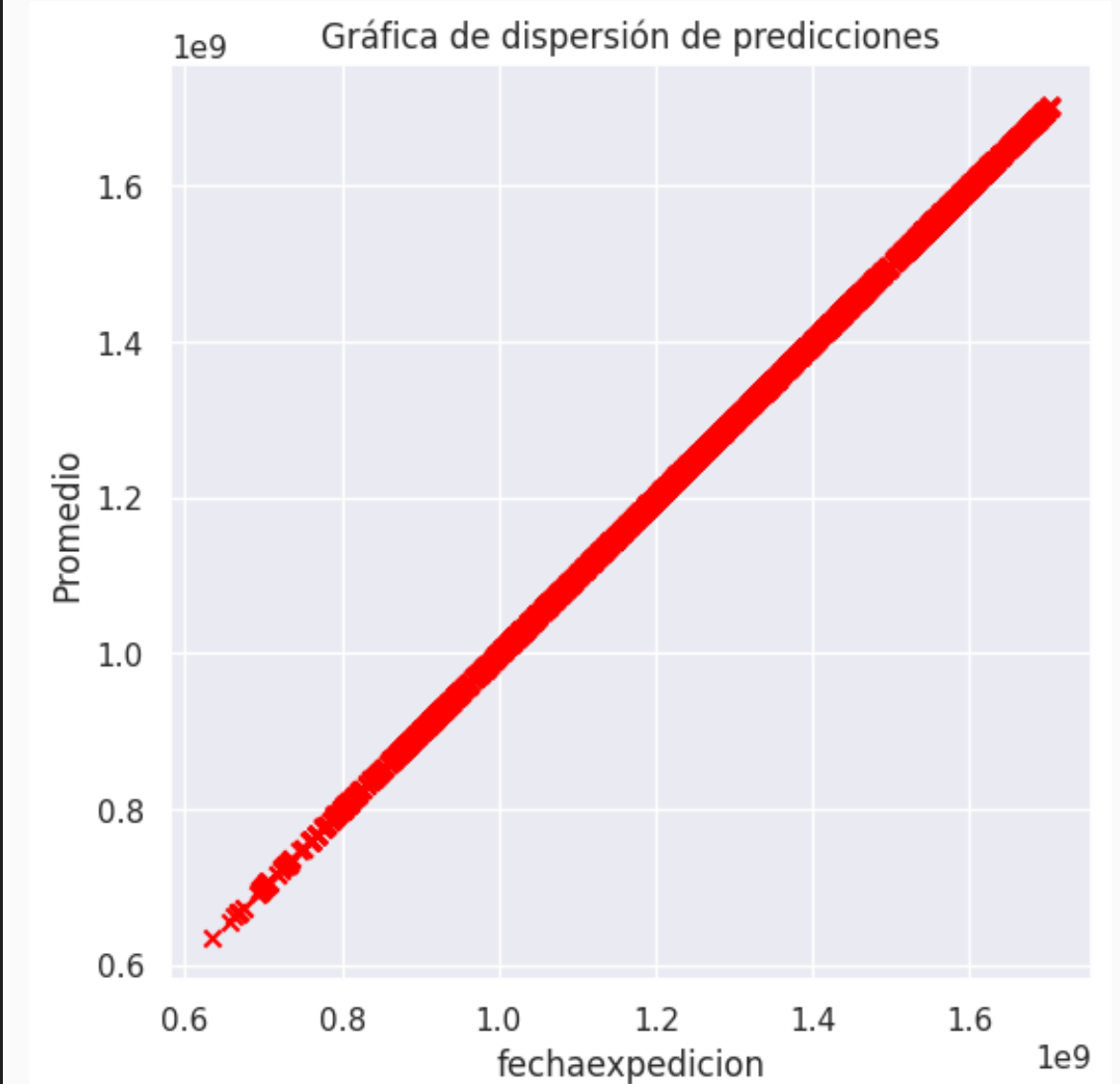


Los valores en el mapa de calor indican la fuerza y la dirección de la relación lineal entre las diferentes variables. El valor cercano a 1 indica una correlación positiva fuerte, mientras que el valor cercano a -1 indica una correlación negativa fuerte. El valor cercano a 0 indica una correlación débil.

```

1 from sklearn.naive_bayes import GaussianNB
2 from sklearn.model_selection import cross_val_score
3 from sklearn.model_selection import train_test_split
4 from sklearn.metrics import accuracy_score
5
6
7 X = datos[['fechaexpedicion']]
8 y = datos[['Promedio']]
9
10 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
11 y_train = np.array(y_train).ravel()
12
13
14 est = GaussianNB()
15
16 est.fit(X_train,y_train)
17 print("%.3f" % accuracy_score(est.predict(X_test), y_test))
18
19 fig = plt.figure(figsize=(6, 6))
20 plt.scatter(X_test['fechaexpedicion'], X_test['fechaexpedicion'], c="red" , cmap='coolwarm', marker='x')
21 plt.xlabel('fechaexpedicion')
22 plt.ylabel('Promedio')
23 plt.title('Gráfica de dispersión de predicciones')
24 plt.show()

```



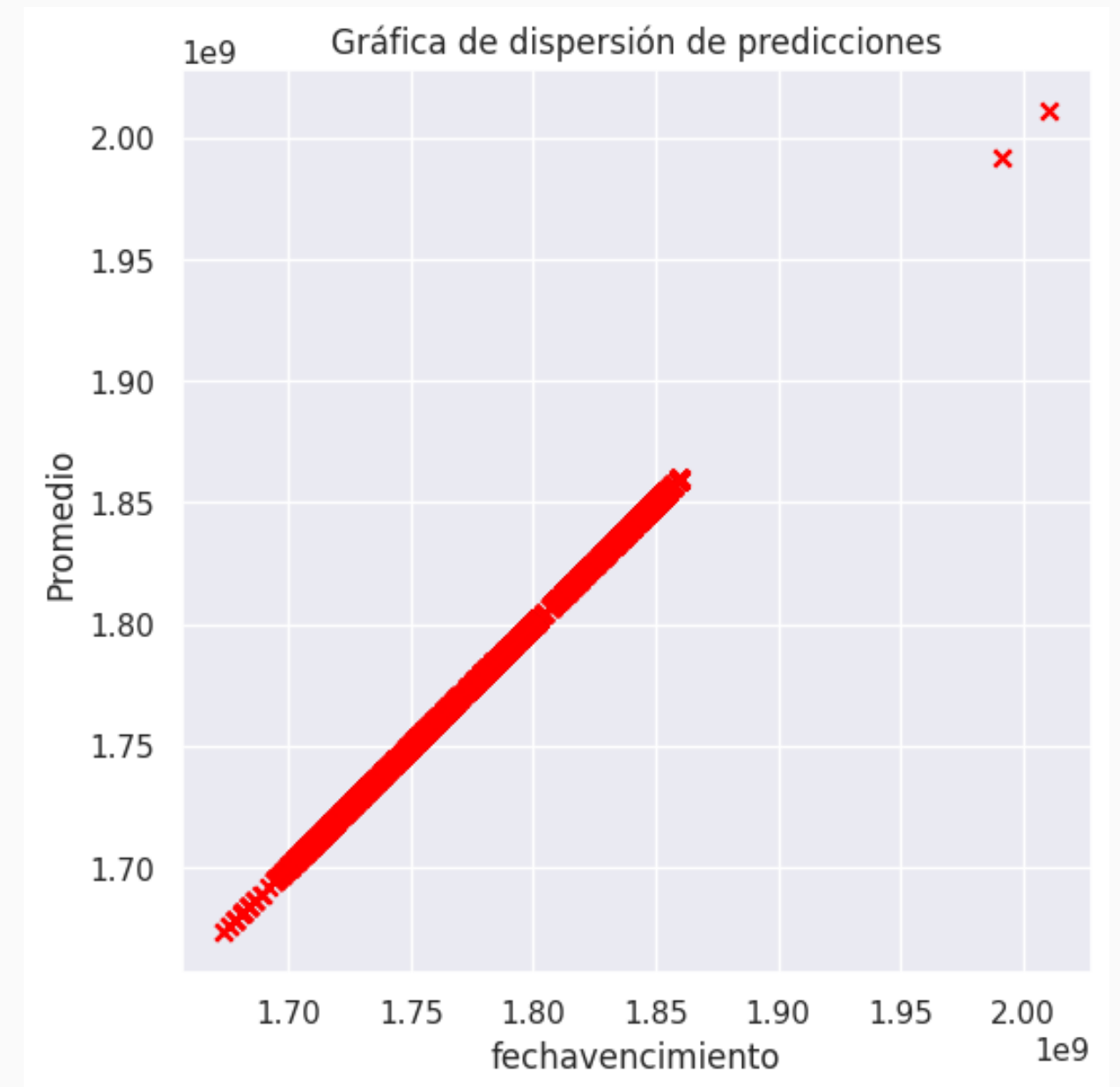
Este código utiliza un Regresor Gaussiano para predecir 'Promedio' basado en la característica 'fechaexpedicion'. La precisión del modelo se imprime y se visualizan las predicciones en un gráfico de dispersión

```

1 from sklearn.naive_bayes import GaussianNB
2 from sklearn.model_selection import cross_val_score
3 from sklearn.model_selection import train_test_split
4 from sklearn.metrics import accuracy_score
5
6
7 X = datos[['fechavencimiento']]
8 y = datos[['Promedio']]
9
10 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
11 y_train = np.array(y_train).ravel()
12
13
14 est = GaussianNB()
15
16 est.fit(X_train,y_train)
17 print("%.3f" % accuracy_score(est.predict(X_test), y_test))
18
19 fig = plt.figure(figsize=(6, 6))
20 plt.scatter(X_test['fechavencimiento'], X_test['fechavencimiento'], c='red', cmap='coolwarm', marker='x')
21 plt.xlabel('fechavencimiento')
22 plt.ylabel('Promedio')
23 plt.title('Gráfica de dispersión de predicciones')
24 plt.show()

```

- Cada punto en la gráfica representa una instancia del conjunto de prueba.
- El color de los puntos (marcadores 'x') indica las predicciones del modelo. Los colores más cálidos o más fríos podrían representar diferentes valores predichos.

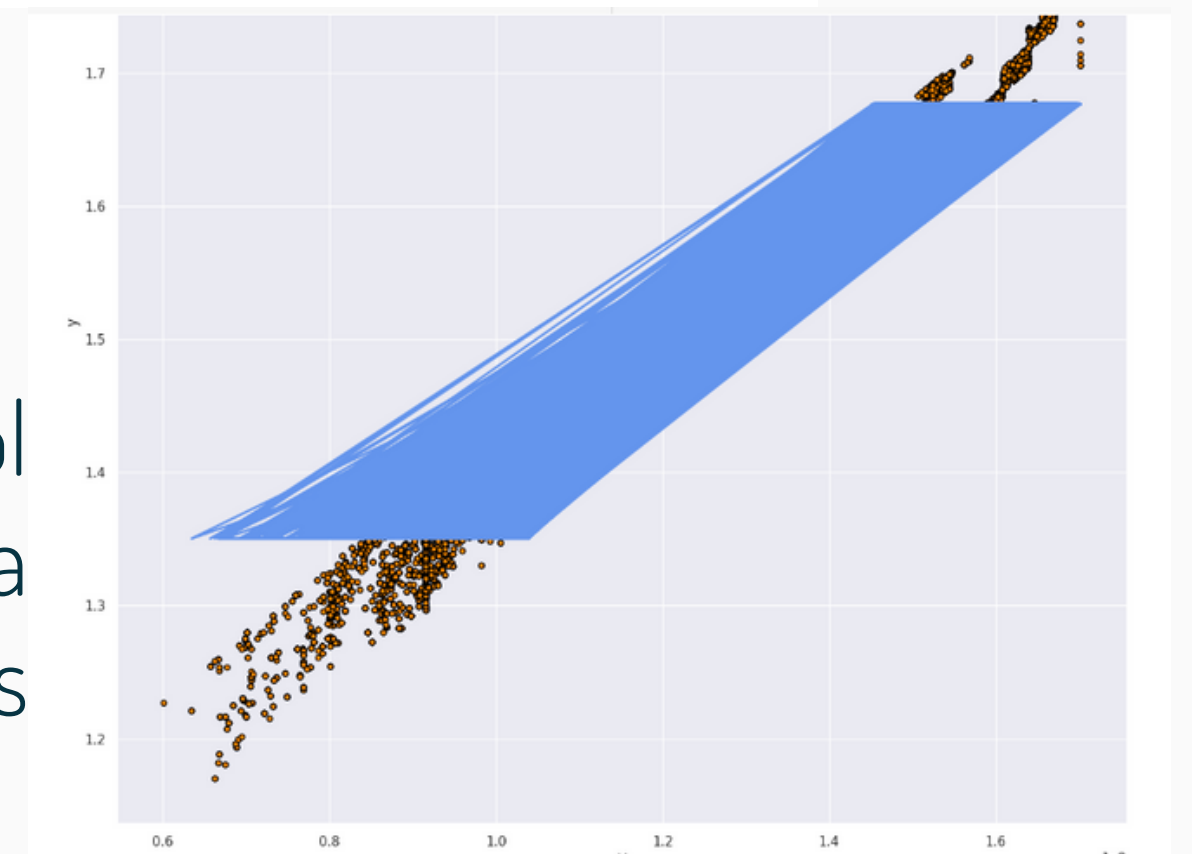
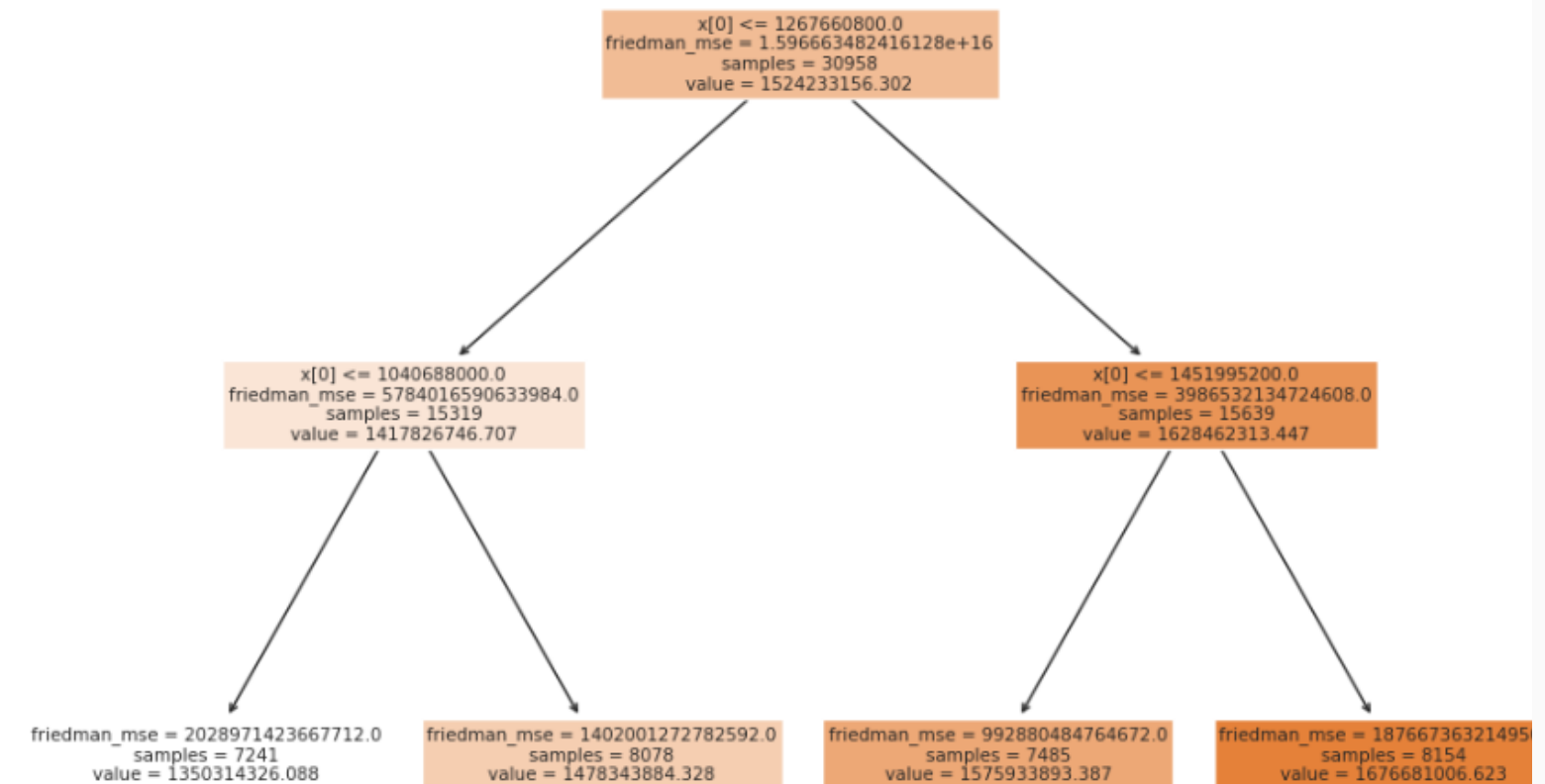



```

1 from sklearn.model_selection import train_test_split
2 from sklearn.tree import DecisionTreeRegressor
3 from sklearn.metrics import mean_squared_error
4 from sklearn.metrics import mean_squared_error, mean_absolute_error
5
6 X = datos[['fechaexpedicion']]
7 y = datos[['Promedio']]
8
9 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
10
11 max_depth = 2
12 criterion = 'friedman_mse'
13 clf = DecisionTreeRegressor(max_depth=max_depth, criterion=criterion)
14
15 clf.fit(X_train, y_train)
16
17 y_pred = clf.predict(X_test)
18
19 mse = mean_squared_error(y_test, y_pred)
20 print("MSE :", mse)
21
22 rmse = np.sqrt(mse)
23 print("RMSE:", rmse)
24
25 y_pred = clf.predict(X_test)
26 mae = mean_absolute_error(y_test, y_pred)
27 print("MAE :", mae)
28
29 from sklearn.tree import plot_tree
30 import matplotlib.pyplot as plt
31
32 plt.figure(figsize=(12, 8))
33 plot_tree(clf, filled=True)
34 plt.show()
35
36 plt.figure()
37 plt.scatter(X, y, s=20, edgecolor="black", c="darkorange", label="data")
38 plt.plot(X_test, y_pred, color="cornflowerblue", linewidth=2, label="regression")
39 plt.xlabel("X")
40 plt.ylabel("y")
41 plt.title("Decision Tree Regression")
42 plt.legend()

```

MSE : 1566960653612978.2
RMSE: 39584853.84099553
MAE : 31748084.99130711



Este código realiza un análisis de regresión utilizando un árbol de decisión, evalúa su rendimiento y proporciona visualizaciones para entender mejor tanto el modelo como sus decisiones.

```

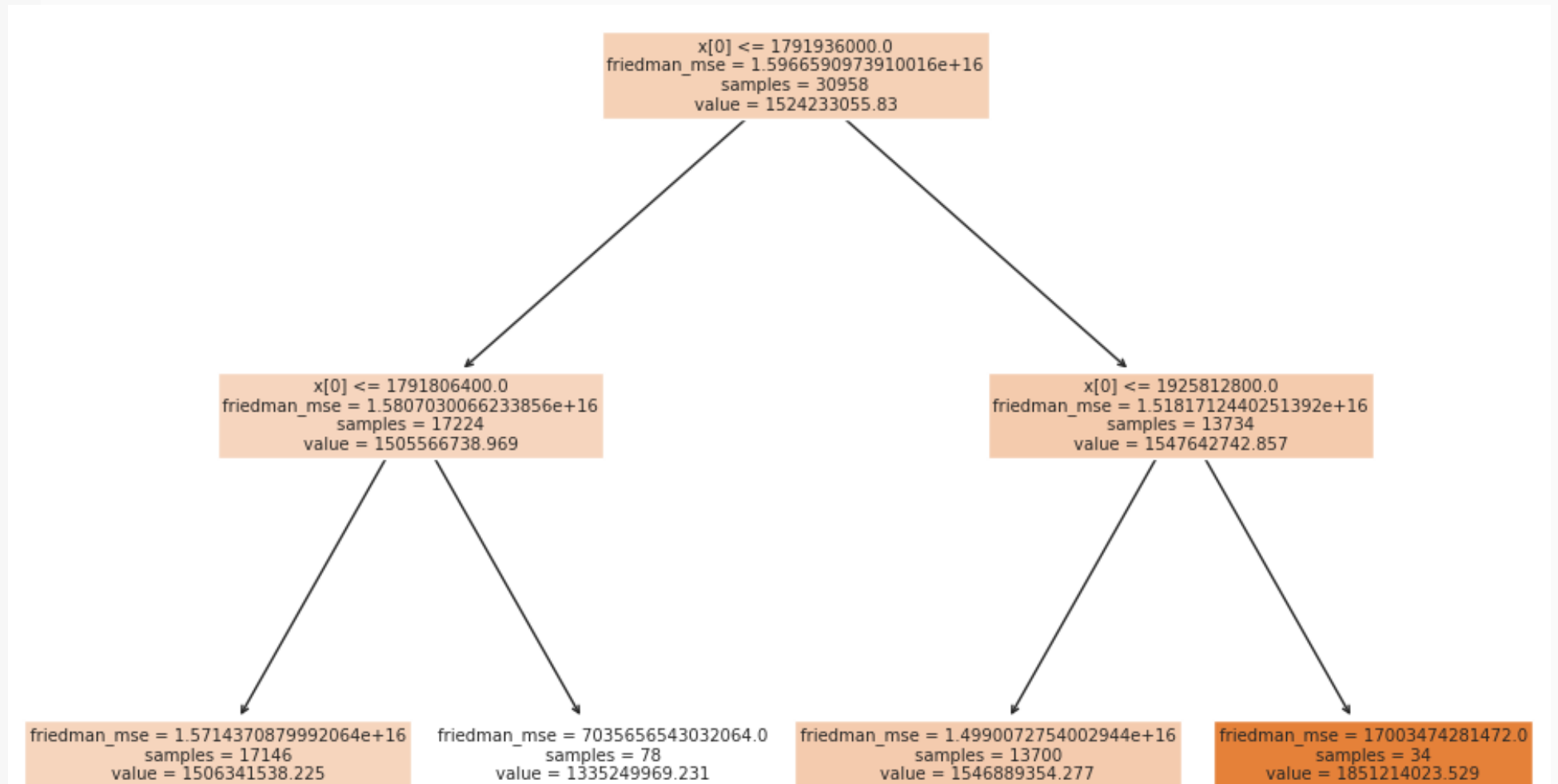
1 from sklearn.model_selection import train_test_split
2 from sklearn.tree import DecisionTreeRegressor
3 from sklearn.metrics import mean_squared_error
4 from sklearn.metrics import mean_squared_error, mean_absolute_error
5
6 X = datos[['fechavencimiento']]
7 y = datos[['Promedio']]
8
9 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
10
11 max_depth = 2
12 criterion = 'friedman_mse'
13 clf = DecisionTreeRegressor(max_depth=max_depth, criterion=criterion)
14
15 clf.fit(X_train, y_train)
16
17 y_pred = clf.predict(X_test)
18
19 mse = mean_squared_error(y_test, y_pred)
20 print("MSE :", mse)
21
22 rmse = np.sqrt(mse)
23 print("RMSE: ", rmse)
24
25 y_pred = clf.predict(X_test)
26 mae = mean_absolute_error(y_test, y_pred)
27 print("MAE :", mae)
28
29 from sklearn.tree import plot_tree
30 import matplotlib.pyplot as plt
31
32 plt.figure(figsize=(12, 8))
33 plot_tree(clf, filled=True)
34 plt.show()

```

```

MSE : 1.5508963703328376e+16
RMSE: 124534989.87565051
MAE : 105395147.13827674

```

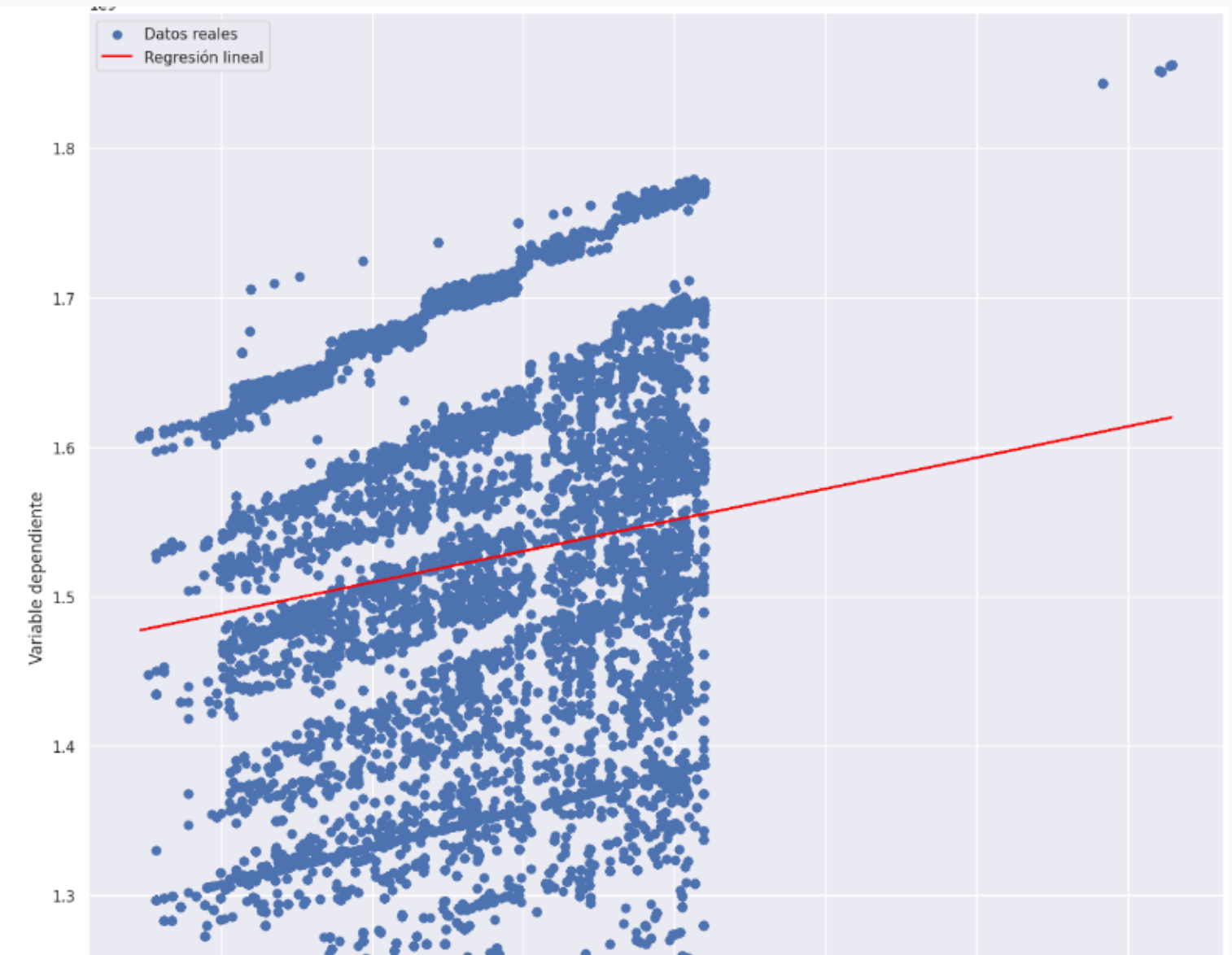


Este código realiza un análisis de regresión utilizando un árbol de decisión, evalúa su rendimiento y proporciona visualizaciones para entender mejor tanto el modelo como sus decisiones específicas en relación con la fecha de vencimiento.

```

1 import numpy as np
2 from sklearn.linear_model import LinearRegression
3
4 X = datos[['fechavencimiento']]
5 y = datos[['Promedio']]
6
7 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
8
9 model = LinearRegression()
10 model.fit(X, y)
11 y_pred = model.predict(X)
12 slope = model.coef_[0]
13 intercept = model.intercept_
14 r_squared = model.score(X, y)
15
16 plt.scatter(X, y, label='Datos reales')
17 plt.plot(X, y_pred, color='red', label='Regresión lineal')
18 plt.legend()
19 plt.xlabel('Variable independiente')
20 plt.ylabel('Variable dependiente')
21 plt.show()

```

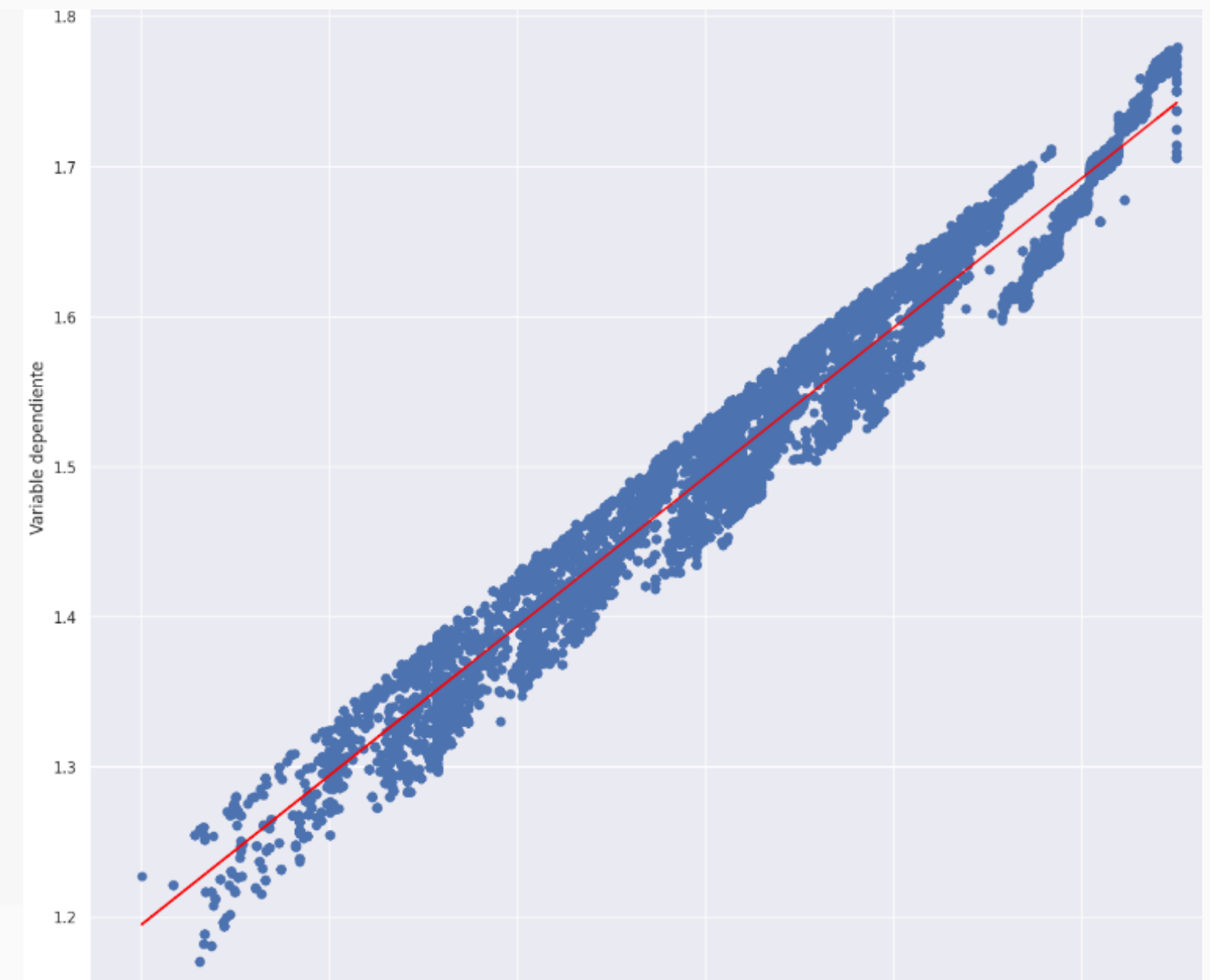


Este código utiliza un modelo de regresión lineal simple para modelar la relación entre la fecha de vencimiento y el promedio. Calcula algunos parámetros clave y visualiza la regresión lineal resultante. Este enfoque es útil para determinar si la relación entre las variables es lineal.

```

1 import numpy as np
2 from sklearn.linear_model import LinearRegression
3
4 X = datos[['fechaexpedicion']]
5 y = datos[['Promedio']]
6
7 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
8
9 model = LinearRegression()
10 model.fit(X, y)
11 y_pred = model.predict(X)
12 slope = model.coef_[0]
13 intercept = model.intercept_
14 r_squared = model.score(X, y)
15
16 plt.scatter(X, y, label='Datos reales')
17 plt.plot(X, y_pred, color='red', label='Regresión lineal')
18 plt.legend()
19 plt.xlabel('Variable independiente')
20 plt.ylabel('Variable dependiente')
21 plt.show()

```

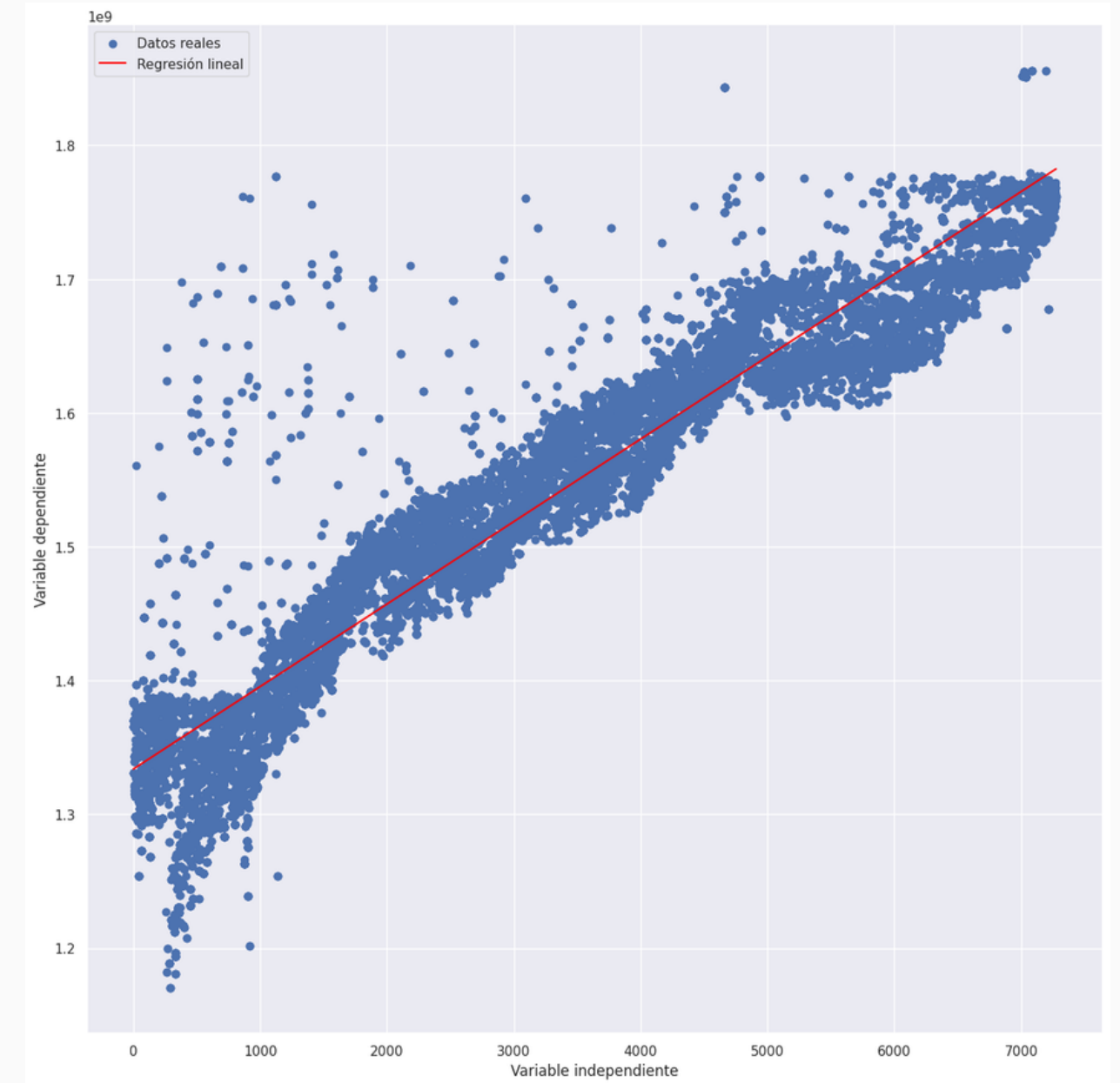


Este código utiliza un modelo de regresión lineal simple para modelar la relación entre la fecha de vencimiento y el promedio. Calcula algunos parámetros clave y visualiza la regresión lineal resultante. Este enfoque es útil para determinar si la relación entre las variables es lineal.


```

1 import numpy as np
2 from sklearn.linear_model import LinearRegression
3
4 X = datos[['producto']]
5 y = datos[['Promedio']]
6
7 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
8
9 model = LinearRegression()
10 model.fit(X, y)
11 y_pred = model.predict(X)
12 slope = model.coef_[0]
13 intercept = model.intercept_
14 r_squared = model.score(X, y)
15
16 plt.scatter(X, y, label='Datos reales')
17 plt.plot(X, y_pred, color='red', label='Regresión lineal')
18 plt.legend()
19 plt.xlabel('Variable independiente')
20 plt.ylabel('Variable dependiente')
21 plt.show()

```



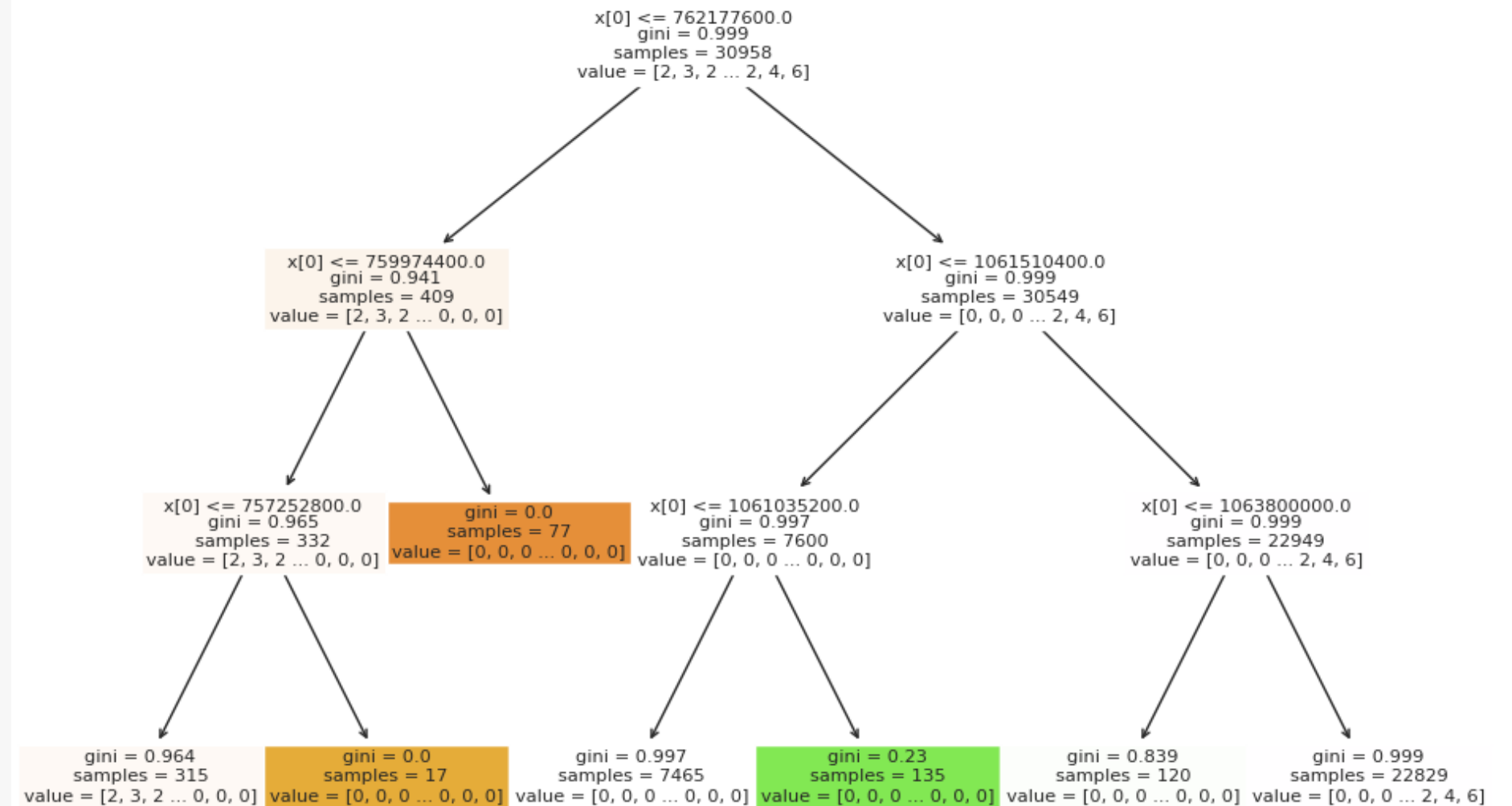
Este código utiliza un modelo de regresión lineal simple para modelar la relación entre el producto y el promedio. Calcula algunos parámetros clave y visualiza la regresión lineal resultante. Este enfoque es útil para determinar si la relación entre las variables es lineal.


```

1 from sklearn.tree import DecisionTreeClassifier
2 from sklearn.model_selection import train_test_split
3 from sklearn.metrics import accuracy_score
4 from sklearn.tree import plot_tree
5 import matplotlib.pyplot as plt
6
7 X = datos[['fechaexpedicion']]
8 y = datos[['Promedio']]
9 np.random.seed(20)
10
11 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
12
13 max_depth = 3
14 criterion = 'gini'
15 clf = DecisionTreeClassifier(max_depth=max_depth, criterion=criterion)
16
17 clf.fit(X_train, y_train)
18
19 y_pred = clf.predict(X_test)
20
21 accuracy = accuracy_score(y_test, y_pred)
22 print(f"Precisión del modelo: {accuracy}")
23
24 plt.figure(figsize=(12, 8))
25 plot_tree(clf, filled=True)
26 plt.show()

```

Precisión del modelo: 0.017054263565891473



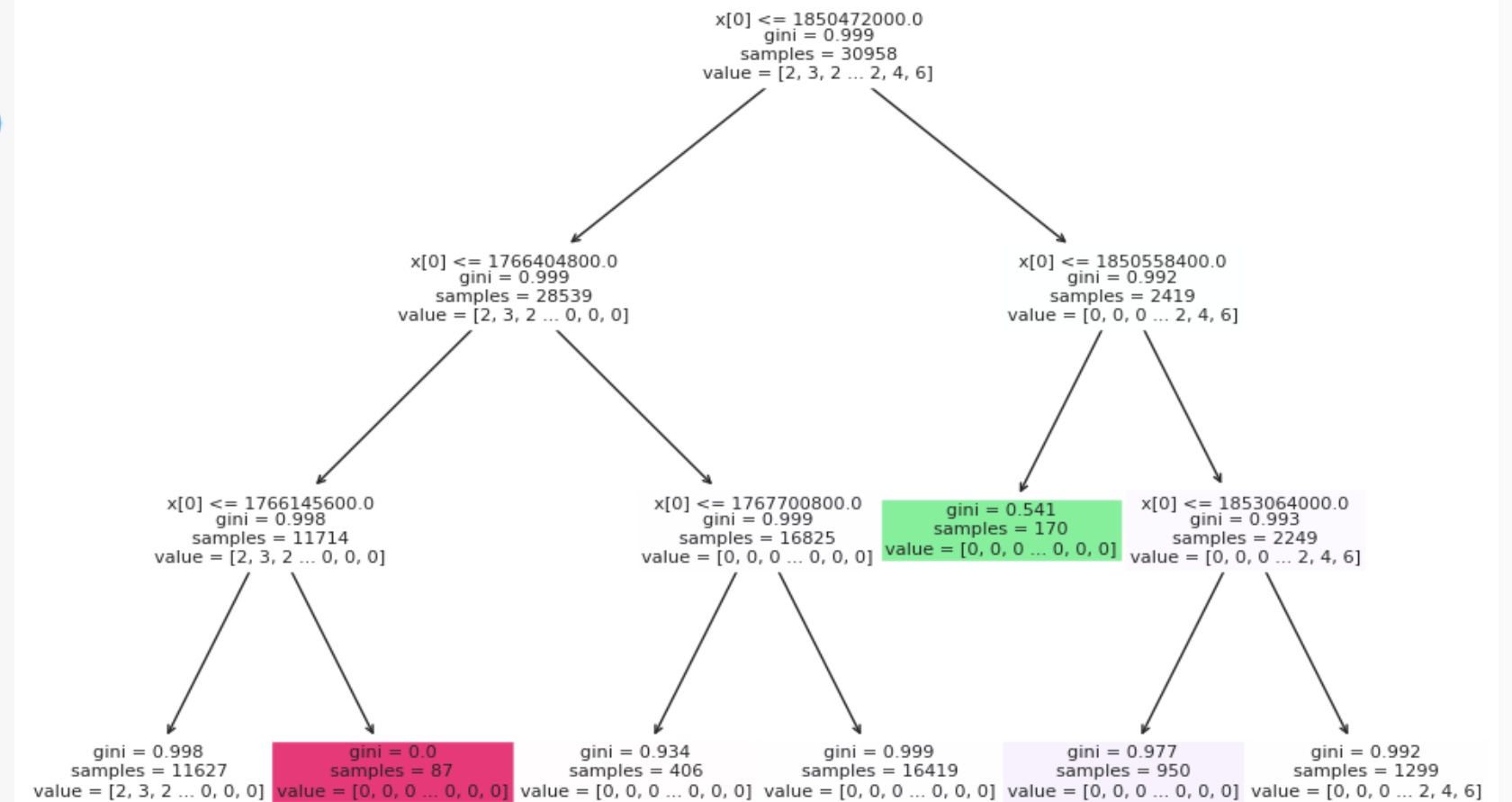
Este código utiliza un Decision Tree Classifier para predecir la variable Promedio basándose en la fecha de expedición. Luego, evalúa la precisión del modelo y visualiza el árbol de decisión para entender cómo toma decisiones.

```

1 from sklearn.tree import DecisionTreeClassifier
2 from sklearn.model_selection import train_test_split
3 from sklearn.metrics import accuracy_score
4
5 X = datos[['fechavencimiento']]
6 y = datos[['Promedio']]
7 np.random.seed(20)
8
9 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
10
11 max_depth = 3
12 criterion = 'gini'
13 clf = DecisionTreeClassifier(max_depth=max_depth, criterion=criterion)
14
15 clf.fit(X_train, y_train)
16
17 y_pred = clf.predict(X_test)
18
19 accuracy = accuracy_score(y_test, y_pred)
20 print(f"Precisión del modelo: {accuracy}")
21
22 from sklearn.tree import plot_tree
23 import matplotlib.pyplot as plt
24
25 plt.figure(figsize=(12, 8))
26 plot_tree(clf, filled=True)
27 plt.show()
28

```

Precisión del modelo: 0.021576227390180877

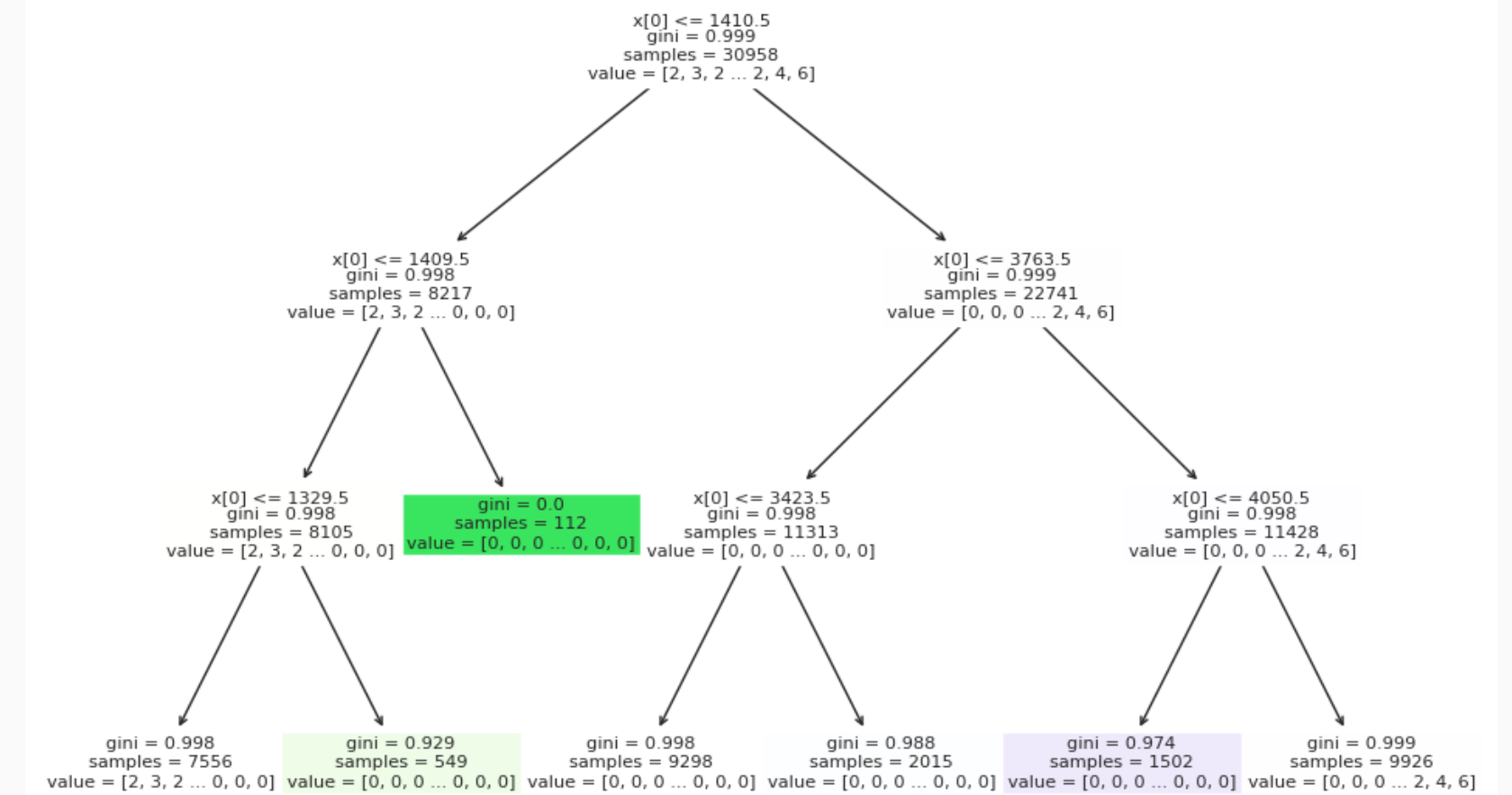


Este código utiliza un Decision Tree Classifier para predecir la variable Promedio basándose en la fecha de vencimiento (fechavencimiento). Evalúa la precisión del modelo y proporciona una visualización del árbol de decisión para entender cómo toma decisiones.

```

1 from sklearn.tree import DecisionTreeClassifier
2 from sklearn.model_selection import train_test_split
3 from sklearn.metrics import accuracy_score
4 from sklearn.tree import plot_tree
5 import matplotlib.pyplot as plt
6
7 X = datos[['producto']]
8 y = datos[['Promedio']]
9 np.random.seed(20)
10
11 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
12
13 max_depth = 3
14 criterion = 'gini'
15 clf = DecisionTreeClassifier(max_depth=max_depth, criterion=criterion)
16
17 clf.fit(X_train, y_train)
18
19 y_pred = clf.predict(X_test)
20
21 accuracy = accuracy_score(y_test, y_pred)
22 print(f"Precisión del modelo: {accuracy}")
23
24 plt.figure(figsize=(12, 8))
25 plot_tree(clf, filled=True)
26 plt.show()

```



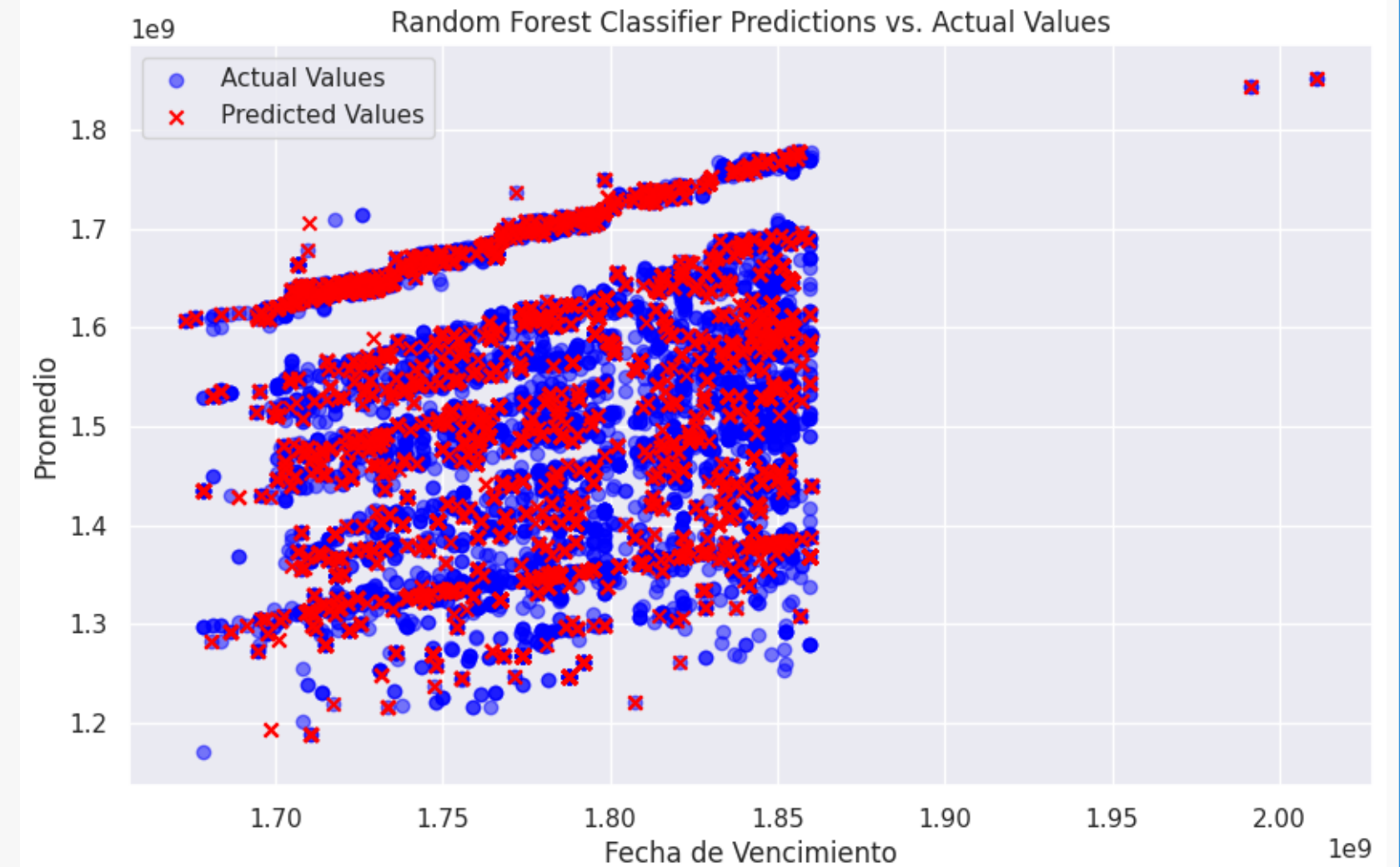
Este código utiliza un Decision Tree Classifier para predecir la variable Promedio basándose en la columna producto. Evalúa la precisión del modelo y proporciona una visualización del árbol de decisión para entender cómo toma decisiones.


```

1 from sklearn.ensemble import RandomForestClassifier
2 from sklearn.model_selection import train_test_split
3 from sklearn.metrics import accuracy_score
4
5 X = datos[['fechavencimiento']]
6 y = datos[['Promedio']]
7
8 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
9
10 y_train = np.array(y_train).ravel()
11 n_estimators = 100
12 max_depth = 400
13 criterion = 'gini'
14 clf = RandomForestClassifier(n_estimators=n_estimators, max_depth=max_depth, criterion=criterion)
15
16 clf.fit(X_train, y_train)
17
18 y_pred = clf.predict(X_test)
19
20
21 accuracy = accuracy_score(y_test, y_pred)
22 print(f"Precisión del modelo: {accuracy}")
23
24 plt.show()
25
26 plt.figure(figsize=(10, 6))
27 plt.scatter(X_test, y_test, color='blue', label='Actual Values', alpha=0.5)
28 plt.scatter(X_test, y_pred, color='red', marker='x', label='Predicted Values')
29 plt.xlabel('Fecha de Vencimiento')
30 plt.ylabel('Promedio')
31 plt.legend()
32 plt.title('Random Forest Classifier Predictions vs. Actual Values')
33 plt.show()

```

Precisión del modelo: 0.3922480620155039



Este código utiliza un Random Forest Classifier o para predecir la variable Promedio basándose en la fecha de vencimiento (fechavencimiento). Evalúa la precisión del modelo y proporciona una visualización de cómo se comparan las predicciones con los valores reales. El bosque aleatorio es una técnica de ensamble que utiliza múltiples árboles de decisión para mejorar la precisión del modelo.

```

1 from sklearn.ensemble import RandomForestClassifier
2 from sklearn.model_selection import train_test_split
3 from sklearn.metrics import accuracy_score
4
5 X = datos[['fechaexpedicion']]
6 y = datos[['Promedio']]
7
8 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
9
10 y_train = np.array(y_train).ravel()
11 n_estimators = 100
12 max_depth = 400
13 criterion = 'gini'
14 clf = RandomForestClassifier(n_estimators=n_estimators, max_depth=max_depth, criterion=criterion)
15
16 clf.fit(X_train, y_train)
17
18 y_pred = clf.predict(X_test)
19
20
21 accuracy = accuracy_score(y_test, y_pred)
22 print(f"Precisión del modelo: {accuracy}")
23
24 plt.show()
25
26 plt.figure(figsize=(10, 6))
27 plt.scatter(X_test, y_test, color='blue', label='Actual Values', alpha=0.5)
28 plt.scatter(X_test, y_pred, color='red', marker='x', label='Predicted Values')
29 plt.xlabel('Fecha de Expedicion')
30 plt.ylabel('Promedio')
31 plt.legend()
32 plt.title('Random Forest Classifier Predictions vs. Actual Values')
33 plt.show()

```

Este código utiliza un Random Forest Classifier para predecir la variable Promedio basándose en la fecha de expedición (fechaexpedicion). Evalúa la precisión del modelo y proporciona una visualización de cómo se comparan las predicciones con los valores reales. El bosque aleatorio es una técnica de ensamble que utiliza múltiples árboles de decisión para mejorar la precisión del modelo.

PCA

```
1 import numpy as np
2 from sklearn.decomposition import PCA
3 import matplotlib.pyplot as plt
4
5 columna_pca = ['producto', 'Promedio']
6
7 data_for_pca = datos[columna_pca]
8
9 mean = np.mean(data_for_pca)
10 std = np.std(data_for_pca)
11 print("mean:", mean)
12 print("std:", std)
13
14 standardized_data = (data_for_pca - mean) / std
```

```
mean: producto    3.086220e+03
Promedio    1.524067e+09
dtype: float64
std: producto    1.944146e+03
Promedio    1.264999e+08
dtype: float64
```

```
/usr/local/lib/python3.10/dist-packages/numpy/core/fromnumeric.py:3430: FutureWarning: In a future version, DataFrame.mean(axis=None) will return a scalar mean over the entire DataFrame. To retain the old behavior, use 'frame.mean(axis=0)' or
return mean(axis=axis, dtype=dtype, out=out, **kwargs)
```

1 standardized_data

	producto	Promedio
0	-1.587443	-1.215856
1	-1.587443	-1.215856
2	-1.587443	-1.215856
8	-1.587443	-1.215856
11	-1.587443	-1.215856
...
114701	2.154046	1.931087
114702	2.154046	1.931087
114703	2.154046	1.931087

Se busca obtener la covarianza de la matriz por medio de la estandarizacion, para despues reducir la data

```
1 covariance_matrix = np.cov(standardized_data, ddof = 1, rowvar = False)
2 covariance_matrix

array([[1.00002584, 0.94747949],
       [0.94747949, 1.00002584]])

1 eigenvalues, eigenvectors = np.linalg.eig(covariance_matrix)
2
3 print("eigenvalues", eigenvalues)
4
5 print("eigenvectors", eigenvectors)

eigenvalues [1.94750533 0.05254635]
eigenvectors [[ 0.70710678 -0.70710678]
 [ 0.70710678  0.70710678]]

1 # np.argsort can only provide lowest to highest; use[::-1] to reverse the list
2 order_of_importance = np.argsort(eigenvalues)[::-1]
3
4 # utilize the sort order to sort eigenvalues and eigenvectors
5 sorted_eigenvalues = eigenvalues[order_of_importance]
6 sorted_eigenvectors = eigenvectors[:,order_of_importance] # sort the columns

1 # use sorted_eigenvalues to ensure the explained variances correspond to the eigenvectors
2 explained_variance = sorted_eigenvalues / np.sum(sorted_eigenvalues)

1 k = 2 # select the number of principal components
2 reduced_data = np.matmul(standardized_data, sorted_eigenvectors[:, :k]) # transform the original data
3
4 print(reduced_data)

      0      1
0    -1.982232  0.262751
1    -1.982232  0.262751
2    -1.982232  0.262751
8     -1.982232  0.262751
11    -1.982232  0.262751
...      ...      ...
114701  2.888625 -0.157656
114702  2.888625 -0.157656
114703  2.888625 -0.157656
114704  2.851801 -0.195208
114705  2.852165 -0.195571

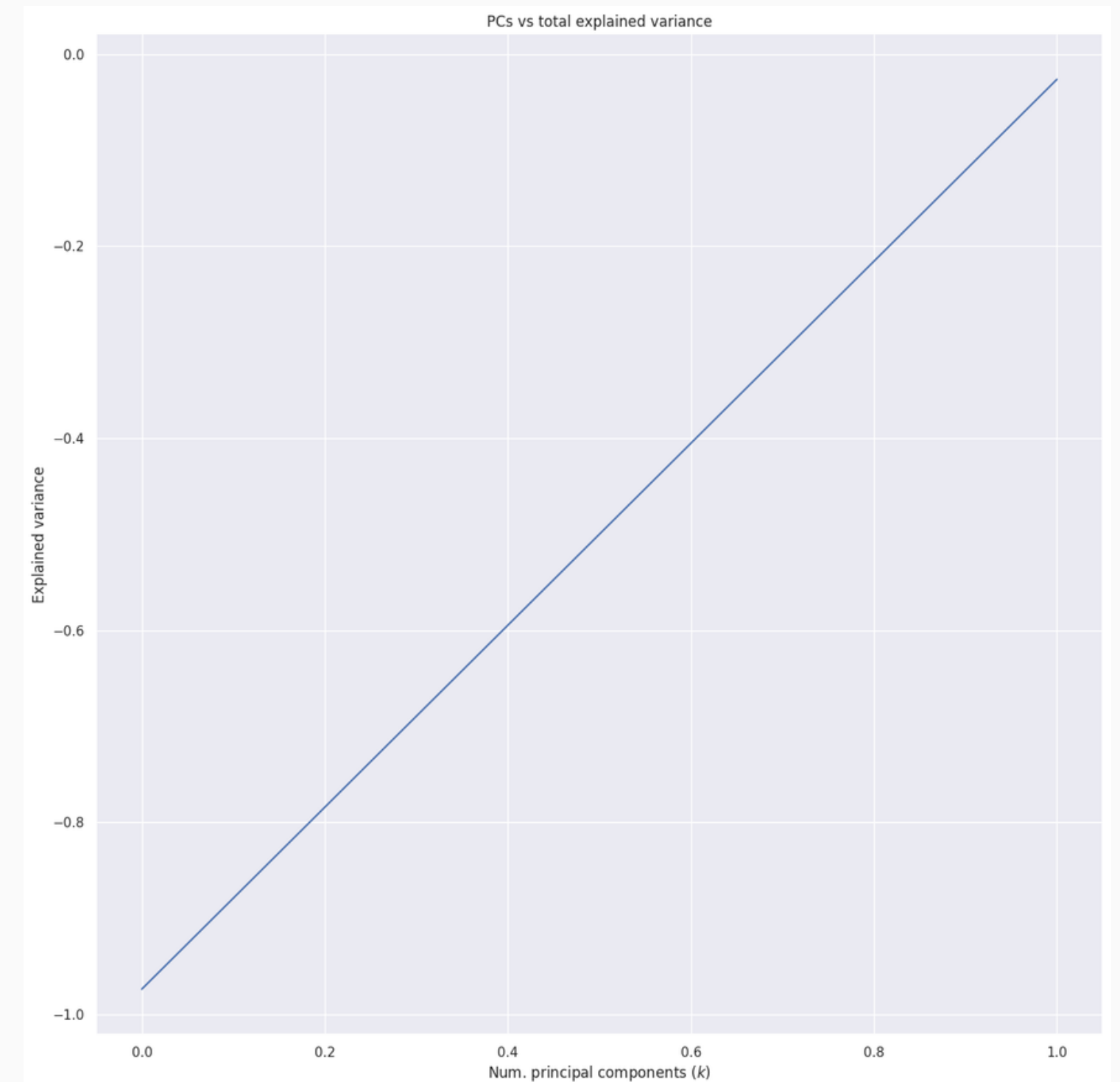
[38698 rows x 2 columns]
```

Aqui se codifica el “Total explained Variance” para poder conocer los n_components

```
[49] 1 total_explained_variance = sum(explained_variance[:k])  
      2  
      3 print(explained_variance)  
      4  
      5 print("total_explained_variance (k=2):", total_explained_variance)
```

```
[0.9737275 0.0262725]  
total_explained_variance (k=2): 1.0
```

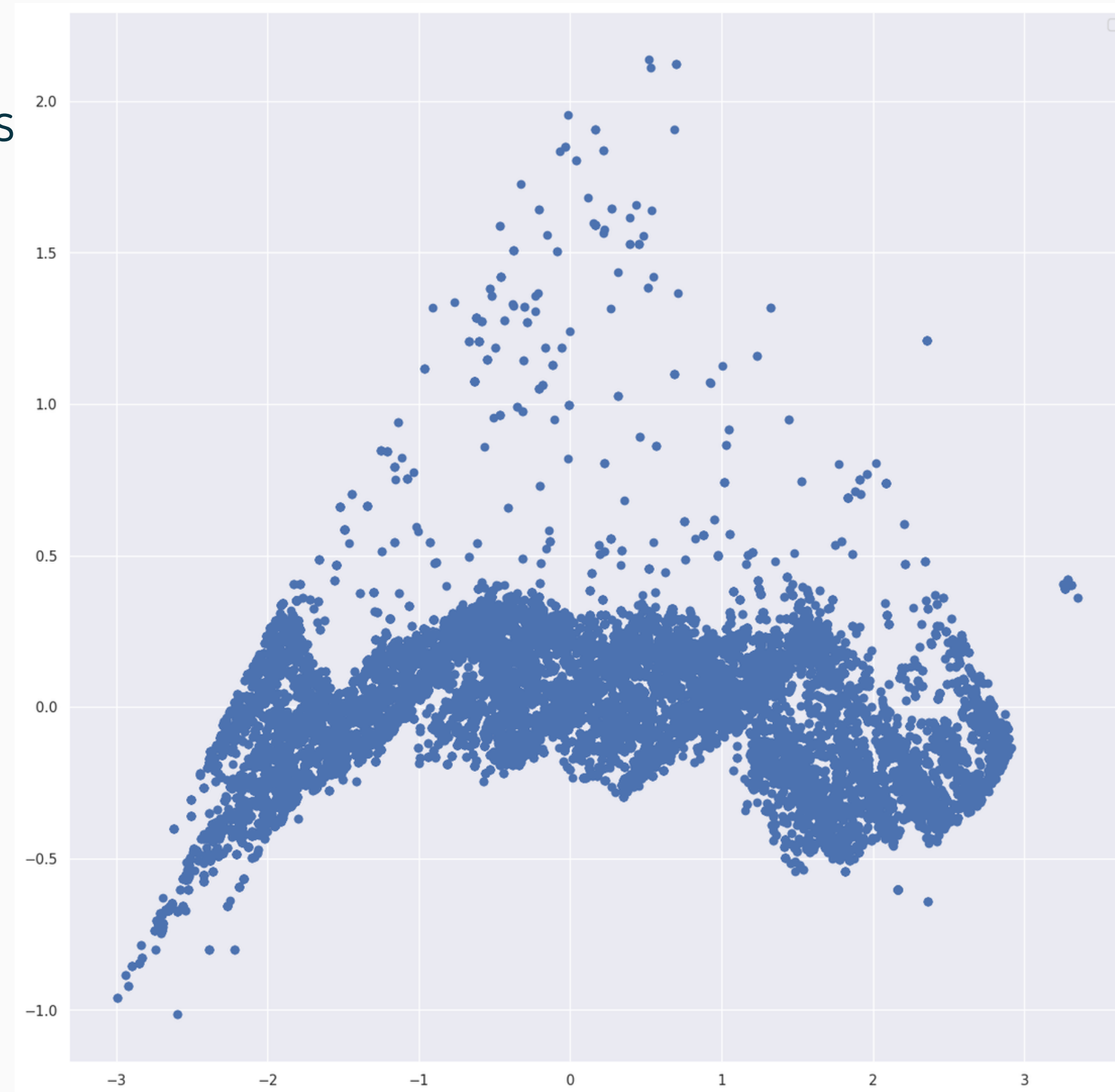
```
[50] 1 plt.plot(-1*explained_variance)  
      2 plt.title("PCs vs total explained variance")  
      3 plt.xlabel("Num. principal components ($k$)")  
      4 plt.ylabel("Explained variance")  
      5 plt.show()
```



Y con los `n_components` identificados lo aplicamos junto a la libreria de PCA para obtener la grafica

```
1 from sklearn.decomposition import PCA
2
3 mypca = PCA(n_components=2)
4 X_pca = mypca.fit_transform(standardized_data)

1 plt.scatter(X_pca[:,0], X_pca[:,1], cmap="rainbow")
2 plt.legend()
3 plt.show()
```



Red Neuronal

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import train_test_split
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.metrics import mean_squared_error
6 import matplotlib.pyplot as plt
7 import tensorflow as tf
8 from tensorflow import keras
9
10 X = datos[['producto']].values
11 y = datos['Promedio'].values
12 scaler = StandardScaler()
13 X_scaled = scaler.fit_transform(X)
14
15 X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.2, random_state=42)
16
17 predicciones = model.predict(X_test)
18
19 mse = mean_squared_error(y_test, predicciones)
20 print(f'Mean Squared Error: {mse}')
21
22 for i in range(10):
23     print(f'Valor real: {y_test[i]}, Predicción: {predicciones[i][0]}')
24
25 plt.scatter(y_test, predicciones)
26 plt.xlabel('Valor Real')
27 plt.ylabel('Predicción')
28 plt.title('Predicciones vs Valores Reales')
29 plt.show()
```

